

23/4/18

Data Mining
Assignment-4Avinash Villineni
A20406697Problem 1.1:-3.1.1) Let $A = \{1, 2, 3, 4\}$, $B = \{2, 3, 5, 7\}$, $C = \{2, 4, 6\}$

$$\text{Sim}(A, B) = \frac{A \cap B}{A \cup B} = \frac{2}{6} = \frac{1}{3}$$

$$\text{Sim}(A, C) = \frac{A \cap C}{A \cup C} = \frac{2}{5}$$

$$\text{Sim}(B, C) = \frac{B \cap C}{B \cup C} = \frac{1}{6}$$

3.2.1)

{ "The", "he", "e m", "mo", "mos", "ost", "st",
"te", "et", "ff" }

or

{ The most effective, most effective way, effective way,
way to represent, to represent documents,
represent document as, document as sets, as sets
for, sets for the, for the purpose }

3.3.3) x $h_1(x)$ $h_2(x)$ $h_3(x)$

a)	0	1	2	2
	1	3	5	1
	2	5	2	0
	3	1	5	5
	4	3	2	4
	5	5	5	3

Document text
model

h_1	h_2	h_3	S_1	S_2	S_3	S_4
1	2	2	0	1	0	1
3	5	1	0	1	0	0
5	2	0	1	0	0	1
1	5	5	0	0	1	0
3	2	4	0	0	1	1
5	5	3	1	0	0	0

Signature Matrix:-

	S_1	S_2	S_3	S_4
h_1	5	1	1	1
h_2	2	2	2	2
h_3	0	1	4	0

b) only h_3 hash function is a true permutation

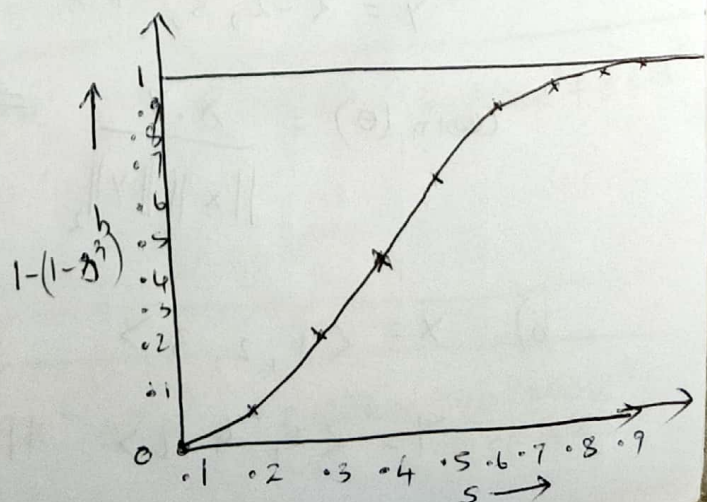
c)

Similarities	S_1-S_2	S_1-S_3	S_1-S_4	S_2-S_3	S_2-S_4	S_3-S_4
col/col	0	0	0.25	0	0.25	0.25
Sig/sig	0.33	0.33	0.67	0.67	0.67	0.67

Jaccard similarities obtained from the Signature matrix using hash functions are not close to the similarities obtained from the Characteristic matrix.

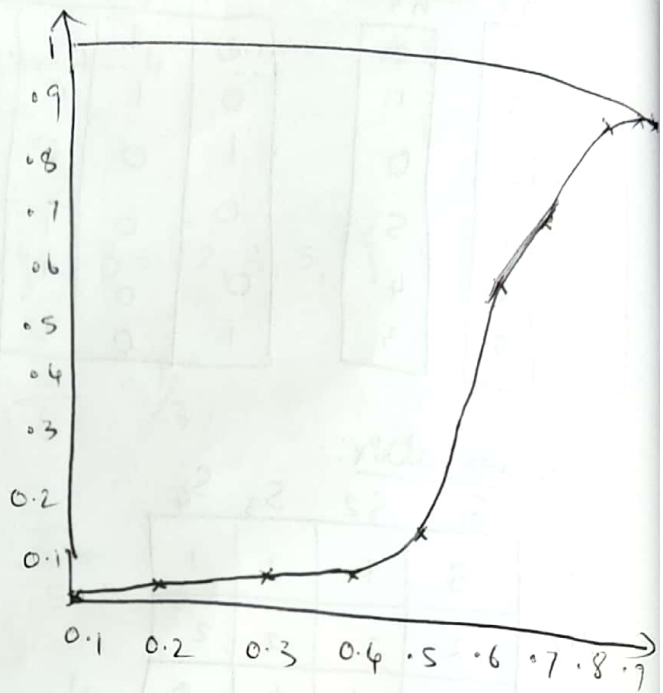
3.4.1) (i) $S = 0.1 \dots 0.9$; $a = 3$ & $b = 10$

S	$1 - (1 - S^a)^b$
0.1	0.01
0.2	0.0772
0.3	0.2394
0.4	0.4839
0.5	0.7369
0.6	0.9123
0.7	0.985
0.8	0.992
0.9	0.999



(11) $a = 6$ & $b = 20$

s	$1 - (1-s)^b$
0.1	0.00002
0.2	0.0013
0.3	0.0145
0.4	0.0788
0.5	0.2702
0.6	0.6154
0.7	0.9182
0.8	0.9977
0.9	1.0



3.5.4)

a) $A = \{1, 2, 3, 4\}$ $B = \{2, 3, 4, 5\}$

$$\text{Sim}(A, B) = \frac{3}{5}$$

b) $A = \{1, 2, 3\}$ $B = \{4, 5, 6\}$

$$\text{Sim}(A, B) = \frac{A \cap B}{A \cup B} = 0$$

3.5.5)

a) $x = \langle 3, -1, 2 \rangle$

$y = \langle -2, 3, 1 \rangle$

$$x \cdot y = -6 - 3 + 2 = -7$$

$$\|x\|_2 = 3.74 \quad \|y\|_2 = 3.74$$

$$\cos(\theta) = \frac{x \cdot y}{\|x\|_2 \|y\|_2} = \frac{-7}{14} = -0.5$$

b) $x = \langle 1, 2, 3 \rangle$

$$x \cdot y = 2 + 8 + 18 = 28$$

$y = \langle 2, 4, 6 \rangle$

$$\|x\|_2 = 3.74 \quad \|y\|_2 = 7.48$$

$$\cos(\theta) = \frac{X \cdot Y}{\|X\| \|Y\|} = \frac{28}{3.74 \times 7.48} = \frac{28}{28} = 1$$

c) $X = \langle 5, 0, -4 \rangle$ $Y = \langle -1, -6, 2 \rangle$ $X \cdot Y = -5 + 0 - 8 = -13$
 $\|X\| = 6.4$ $\|Y\| = 6.4$

$$\cos(\theta) = \frac{-13}{41} = -0.317$$

d) $X = \langle 0, 1, 1, 0, 1, 1 \rangle$ $Y = \langle 0, 0, 1, 0, 0, 0 \rangle$ $X \cdot Y = 1$
 $\|X\| = 2$ $\|Y\| = 1$

$$\cos(\theta) = \frac{1}{2 \times 1} = \frac{1}{2} = 0.5$$

Problem 1.2:-

9.21)

	A	B	C
PS	3.06	2.68	2.92
DS	500	320	640
MS	6	4	6

$$\cos(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$$

$$\cos(A, B) = \frac{8.2008 + 16000d^2 + 24B^2}{\sqrt{9.3636 + 25000d^2 + 36B^2} \sqrt{7.1824 + 102400d^2 + 16B^2}}$$

$$\sqrt{9.3636 + 25000d^2 + 36B^2} \quad \sqrt{7.1824 + 102400d^2 + 16B^2}$$

$$\cos(A, C) = \frac{8.9352 + 32000d^2 + 36B^2}{\sqrt{9.3636 + 25000d^2 + 36B^2} \sqrt{8.5264 + 409600d^2 + 36B^2}}$$

$$\sqrt{9.3636 + 25000d^2 + 36B^2} \quad \sqrt{8.5264 + 409600d^2 + 36B^2}$$

$$\cos(B, C) = \frac{7.8256 + 204800d^2 + 24B^2}{\sqrt{7.1824 + 102400d^2 + 16B^2} \sqrt{8.5264 + 409600d^2 + 36B^2}}$$

$$\sqrt{7.1824 + 102400d^2 + 16B^2} \quad \sqrt{8.5264 + 409600d^2 + 36B^2}$$

$$\cosim(A, B)_{\alpha, \beta=1} = \frac{160032.2}{500.045 \times 320.036} = \frac{160032.2}{160032.5} = 0.99999773$$

$$\cosim(A, C)_{\alpha, \beta=1} = 0.9999953$$

$$\cosim(B, C)_{\alpha, \beta=1} = 0.9999879$$

c) $\alpha = 0.01, \beta = 0.5$

$$\cosim(A, B)_{\alpha=0.01, \beta=0.5} = 0.9908696$$

$$\cosim(A, C)_{\alpha=0.01, \beta=0.5} = 0.9990992$$

$$\cosim(B, C)_{\alpha=0.01, \beta=0.5} = 0.9861894$$

d) $\alpha = \frac{1}{486.66} = 0.0020548$

$$\beta = \frac{1}{5.3} = 0.1875$$

$$\cosim(A, B) = 0.9943904$$

$$\cosim(A, C) = 0.995614$$

$$\cosim(B, C) = 0.9822469$$

9.23)

A	B	C
3.06	2.68	2.92
500	320	640
6	4	6

weight $\Rightarrow A \rightarrow 4, B \rightarrow 2, C \rightarrow 5$

$$\text{avg} = \frac{4+2+5}{3} = 1\frac{1}{3}$$

A	B	C
4	2	5

row avg

$$1\frac{1}{3}$$

Normalized rating

A	B	C
$\frac{1}{3}$	$-\frac{5}{3}$	$\frac{4}{3}$

$$\begin{aligned} \text{b) processor speed} &= 3.06 \times \frac{1}{3} - 2.68 \times \frac{5}{3} + 2.92 \times \frac{4}{3} \\ &= 1.02 - 4.466 + 3.893 \\ &= 0.447 \end{aligned}$$

$$\begin{aligned} \text{Disk size} &= 500 \times \frac{1}{3} - 320 \times \frac{5}{3} + 640 \times \frac{4}{3} \\ &= 166.66 - 533.33 + 853.33 \\ &= 486.667 \end{aligned}$$

$$\begin{aligned} \text{Memory size} &= 6 \times \frac{1}{3} - 4 \times \frac{5}{3} + 6 \times \frac{4}{3} \\ &= 3.333 \end{aligned}$$

~~2.333~~

9.3.1
a)

	a	b	c	d	e	f	g	h
A	4	5		5	1		3	2
B		3	4	3	1	2	1	
C	2		1	3		4	5	3

consider as booleans

$$\text{sim}(A, B) = \frac{4}{8} = \frac{1}{2} ; \text{sim}(A, C) = \frac{4}{8} = \frac{1}{2} = 0.5$$

$$\text{sim}(B, C) = \frac{4}{8} = \frac{1}{2}$$

$$b) \cosim(A, B) = \frac{4}{\sqrt{6} * \sqrt{6}} = 0.6667 ; \cosim(A, C) = \frac{4}{\sqrt{6} * \sqrt{6}} = 0.6667$$

$$\cosim(B, C) = \frac{4}{\sqrt{6} * \sqrt{6}} = 0.6667$$

c)

	a	b	c	d	e	f	g	h
A	1	1	0	1	0	0	1	0
B	0	1	1	1	0	0	0	0
C	0	0	0	1	0	1	1	1

$$J(A, B) = \frac{2}{5} \quad J_D = 1 - \frac{2}{5} = \frac{3}{5}$$

$$J(B, C) = \frac{1}{6} \quad J_D = 1 - \frac{1}{6} = \frac{5}{6}$$

$$J(A, C) = \frac{2}{6} = \frac{1}{3} \quad J_D = 1 - \frac{1}{3} = \frac{2}{3}$$

$$d) \cosim(A, B) = \frac{A \cdot B}{\|A\|_2 \|B\|_2} = \frac{1}{\sqrt{3}} = 0.5777$$

$$\cosim(B, C) = \frac{1}{\sqrt{3} \sqrt{4}} = 0.288$$

$$\cosim(A, C) = \frac{1}{\sqrt{4} \sqrt{4}} = \frac{1}{4} = 0.25$$

e)

$$\text{avg}(A) = \frac{10}{3} \quad \text{avg}(C) = 3$$

$$\text{avg}(B) = \frac{7}{3}$$

	a	b	c	d	e	f	g	h
A	2/3	5/3		5/3	-1/3		-1/3	-4/3
B		2/3	5/3	2/3	-4/3	-1/3	-4/3	
C	-1		-2	0		1	2	0

$$\cos(A, B) = \frac{52/9}{3.6515 \times 2.708} = 0.5843$$

$$\cos(A, C) = -0.11547$$

$$\cos(B, C) = -0.739574$$

Problem 1.3 :-

S.1.1) Transition Matrix

$$\begin{bmatrix} 1/3 & 1/2 & 0 \\ 1/3 & 0 & 1/2 \\ 1/3 & 1/2 & 1/2 \end{bmatrix}$$

$$\lambda = \left[\frac{3}{13}, \frac{4}{13}, \frac{6}{13} \right]^T$$

list

$$\begin{bmatrix} 0.33 \\ 0.33 \\ 0.33 \end{bmatrix}, \begin{bmatrix} 0.277 \\ 0.277 \\ 0.444 \end{bmatrix}, \begin{bmatrix} 0.2314 \\ 0.3148 \\ 0.4537 \end{bmatrix}, \begin{bmatrix} 0.2345 \\ 0.3040 \\ 0.4614 \end{bmatrix},$$

$$\begin{bmatrix} 0.2301 \\ 0.3088 \\ 0.4609 \end{bmatrix} \dots \dots \begin{bmatrix} 0.2307 \\ 0.3076 \\ 0.4615 \end{bmatrix}$$

S.1.2)

$$V' = \beta M V + (1-\beta) e/m$$

$$= \begin{bmatrix} 4/15 & 2/5 & 0 \\ 4/15 & 0 & 2/5 \\ 4/15 & 2/5 & 2/5 \end{bmatrix} V + \begin{bmatrix} 1/15 \\ 1/15 \\ 1/15 \end{bmatrix}$$

$$V = \left[\frac{7}{27}, \frac{25}{81}, \frac{35}{81} \right]^T$$

list:-

$$\begin{bmatrix} 0.33 \\ 0.33 \\ 0.33 \end{bmatrix}, \begin{bmatrix} 0.2888 \\ 0.2888 \\ 0.4222 \end{bmatrix}, \begin{bmatrix} 0.2592 \\ 0.3125 \\ 0.4281 \end{bmatrix}, \begin{bmatrix} 0.2608 \\ 0.3070 \\ 0.4320 \end{bmatrix},$$

$$\begin{bmatrix} 0.2590 \\ 0.3090 \\ 0.4318 \end{bmatrix}, \dots, \begin{bmatrix} 0.2592 \\ 0.3086 \\ 0.4320 \end{bmatrix}$$

S.1.b)

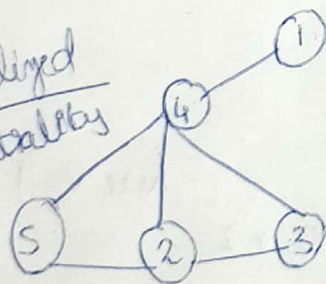
There exist only one head node with self-direction and pageRank of 1. PageRank of the remaining nodes will be $\frac{1}{2}$.

Problem 4 Centrality Measures

$$C_D(V) = \deg(V)$$

a)

Normalized
Degree Centrality



$$C_D^*(V) = \frac{1}{n-1} C_D(V)$$

$$C_D^*(1) = \frac{1}{4} = 0.25$$

$$C_D^*(2) = \frac{3}{4} = 0.75$$

$$C_D^*(5) = \frac{2}{4} = 0.5$$

$$C_D^*(3) = \frac{2}{4} = 0.5$$

$$C_D^*(4) = \frac{4}{4} = 1$$

(ii) Closeness Centrality

$$C_c^*(v) = (n-1) \frac{1}{\sum_j d(v,j)}$$

$$C_c^*(1) = 4 * \frac{1}{7} = \frac{4}{7} = 0.571$$

$$C_c^*(2) = 4 * \frac{1}{5} = \frac{4}{5} = 0.8$$

$$C_c^*(3) = 4 * \frac{1}{6} = \frac{4}{6} = 0.667$$

$$C_c^*(4) = 4 * \frac{1}{4} = 1$$

$$C_c^*(5) = 4 * \frac{1}{6} = 0.667$$

(iii) Betweenness Centrality

A B C D E
1 2 3 4 5

Betweenness of 1:-

$$2 \rightarrow 3 \quad 0/1$$

$$2 \rightarrow 4 \quad 0/1$$

$$2 \rightarrow 5 \quad 0/1$$

$$3 \rightarrow 4 \quad 0/1$$

$$3 \rightarrow 5 \quad 0/2$$

$$4 \rightarrow 5 \quad 0/1$$

Total = 0

here $v=1$

$$C_B(v) = 0 * 2$$

$$C_B^*(v) = \frac{0 * 2}{2 * (4 \log 2)} = \frac{0 * 2}{12} = 0$$

Betweenness of 2:-

$$1 \rightarrow 3 \quad 0/1$$

$$1 \rightarrow 4 \quad 0/1$$

$$1 \rightarrow 5 \quad 0/1$$

$$3 \rightarrow 4 \quad 0/1$$

$$3 \rightarrow 5 \quad 1/2$$

$$4 \rightarrow 5 \quad 0/1$$

$$\text{Total} = \frac{1}{2}$$

$$CB(2) = \frac{1}{2} \times 2 = 1$$

$$CB^*(2) = \frac{1}{12} = 0.0833$$

Betweenness of 3:-

$$1 \rightarrow 2 \quad 0/1$$

$$1 \rightarrow 4 \quad 0/1$$

$$1 \rightarrow 5 \quad 0/1$$

$$2 \rightarrow 4 \quad 0/1$$

$$2 \rightarrow 5 \quad 0/1$$

$$4 \rightarrow 5 \quad 0/1$$

Total = 0

$$CB^*(3) = 0$$

Betweenness of 4

$$1 \rightarrow 2 \quad 1/1$$

$$1 \rightarrow 3 \quad 1/1$$

$$1 \rightarrow 5 \quad 1/1$$

$$2 \rightarrow 3 \quad 0/1$$

$$2 \rightarrow 5 \quad 0/1$$

$$3 \rightarrow 5 \quad 1/2$$

Total = 3.5

$$CB(4) = 3.5 \times 2 = 7$$

$$CB^*(4) = \frac{7}{12} = 0.5833$$

Betweenness of 5:-

$$1 \rightarrow 2 \quad 0/1$$

$$1 \rightarrow 3 \quad 0/1$$

$$1 \rightarrow 4 \quad 0/1$$

$$2 \rightarrow 3 \quad 0/1$$

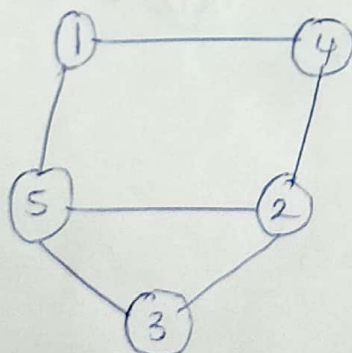
$$2 \rightarrow 4 \quad 0/1$$

$$3 \rightarrow 4 \quad 0/1$$

Total = 0

$$CB^*(5) = 0$$

b)



Degree Centrality:-

$$CD^*(1) = \frac{2}{4} = 0.5$$

$$CD^*(2) = \frac{3}{4} = 0.75$$

$$CD^*(3) = \frac{2}{4} = 0.5$$

$$CD^*(4) = \frac{2}{4} = 0.5$$

$$CD^*(5) = \frac{3}{4} = 0.75$$

(ii) Closeness Centrality

$$C_c^*(v) = (n-1) \cdot \frac{1}{\sum_j d(v, j)}$$

$$= \underline{\underline{\frac{4}{5}}}$$

$$C_c^*(1) = \frac{4}{6} = 0.667$$

$$C_c^*(2) = \frac{4}{5} = 0.8$$

$$C_c^*(3) = \frac{4}{6} = 0.667$$

$$C_c^*(4) = \frac{4}{6} = 0.667$$

$$C_c^*(5) = 0.8 \Rightarrow 4/5$$

(iii) Betweenness Centrality

1 2 3 4 5

Betweenness of 1:-

$$2 \rightarrow 3 \quad 0/1 \quad 3 \rightarrow 4 \quad 0/1$$

$$2 \rightarrow 4 \quad 0/1 \quad 3 \rightarrow 5 \quad 0/1$$

$$2 \rightarrow 5 \quad 0/1 \quad 4 \rightarrow 5 \quad 1/2$$

$$\text{Total} = 0.5$$

$$C_B(1) = 2 \times 0.5 = 1$$

$$C_B^*(1) = \frac{1}{12} = 0.0833$$

Betweenness of 2:-

$$1 \rightarrow 3 \quad 0/1 \quad 3 \rightarrow 4 \quad 1/1$$

$$1 \rightarrow 4 \quad 0/1 \quad 3 \rightarrow 5 \quad 0/1$$

$$1 \rightarrow 5 \quad 0/1 \quad 4 \rightarrow 5 \quad 1/2$$

$$\text{Total} = 1.5$$

$$C_B(2) = 3$$

$$C_B^*(2) = \frac{3}{12} = 0.25$$

Betweenness of 3:-

$$1 \rightarrow 2 \quad 0/2 \quad 2 \rightarrow 4 \quad 0/1$$

$$1 \rightarrow 4 \quad 0/1 \quad 2 \rightarrow 5 \quad 0/1$$

$$1 \rightarrow 5 \quad 0/1 \quad 4 \rightarrow 5 \quad 0/2$$

$$\text{Total} = 0$$

$$C_B^*(3) = 0$$

Betweenness of 4:-

$1 \rightarrow 2$	$\frac{1}{2}$	$2 \rightarrow 3$	$0/1$
$1 \rightarrow 3$	$0/1$	$2 \rightarrow 5$	$0/1$
$1 \rightarrow 5$	$0/1$	$3 \rightarrow 5$	$0/1$

Total = 0.5

$CB(4) = 1$

$$CB^*(4) = \frac{1}{12} = 0.0833$$

Betweenness of 5:-

$1 \rightarrow 2$	$\frac{1}{2}$	$2 \rightarrow 3$	$0/1$
$1 \rightarrow 3$	$\frac{1}{1}$	$2 \rightarrow 4$	$0/1$
$1 \rightarrow 4$	$0/1$	$3 \rightarrow 4$	$0/1$

Total = 1.5

$CB(5) = 3$

$$CB^*(5) = \frac{3}{12} = 0.25$$